



# Potential of I/O-Aware Workflows in Climate and Weather

Julian M. Kunkel<sup>1</sup> , Luciana R. Pedro<sup>1</sup> , *Please add yourself...*

© The Authors 2020. This paper is published with open access at SuperFri.org

The efficient, convenient, and robust execution of data-driven workflows and enhanced data management are key for productivity in scientific computing. Traditionally, in HPC, the concerns of storage and computing are separated and optimised independently from each other and the needs of the end-to-end user.

As climate and weather workflows become increasingly complex and blended beyond data centres while, at the same time, storage hierarchies become deeper, the community investigates ways to reorganise storage access to utilise such systems fully.

The key contributions of this paper are: 1) we sketch the vision of an integrated data-driven approach and discuss the challenges and implications of this strategy, and 2) we illustrate architecture and roadmap that allows the seamless integration into current climate and weather workflows. The tools employed here to achieve an extended workflow are Cylc, XIOS, DDN IME, and ESDM.

We believe workflows composed of data, computing, and communication-intensive tasks should drive the interfaces and hardware configurations to best support the programming models. The changes proposed here increase the opportunity of implementations for smarter scheduling of computing and storage in heterogeneous storage systems.

*Keywords: workflow, climate, weather, heterogeneous storage, data-driven.*

## Introduction

High-Performance Computing (HPC) harnesses the fastest available hardware components to enable the execution of tightly coupled applications from science and industry. Typical use-cases include numerical simulation of physical systems and analysis of large-scale observational data. In the domain of climate and weather, there is a considerable demand for the orchestration of ensembles of simulation models and the generation of data products. A weather forecast service such as the operational weather forecast workflow in Met Office writes around 200TB every day and reads around 600 TB. In total, at the Met Office, on average 1.5 PBs and 14 PBs per day are written and read, respectively, for all climate and weather forecasts across all HPC clusters.

Based on their needs, the HPC community has developed a software ecosystem that supports climate and weather scientists to execute their large-scale workflows. While the current advances correspond to a big leap forward, many processes still require experts. For example, porting a workflow from one system to another still requires adjusting runtime parameters of applications and deciding on how data is managed.

Since performance is of crucial importance to large-scale workflows, careful attention must be paid to exploit the system characteristics of the target supercomputer. A specific supercomputer typically requires substantial changes to the workflow to tailor it to the particular machine and obtain the best performance. For instance, a data-driven workflow may benefit from using a heterogeneous set of computing and storage technology at the same time.

Knowing the capabilities, interfaces, and performance characteristics of individual components are mandatory to make the best use of them. As the complexity of systems expands and alternative storage and computing technologies provide unique characteristics, it becomes increasingly difficult even for experts to manually optimise resource usage in workflows. In many cases, modifications are not performed because: 1) They are labour intense: any change to the

---

<sup>1</sup>University of Reading, Reading, United Kingdom

workflow requires careful validation which may not pay off for small scale runs; 2) Users are not aware of the potential of the complex system; and, 3) because it is a one-time explorative workflow. In this paper, we illustrate how knowing the Input/Output (I/O) characteristics of workflow tasks and overall experimental design helps to optimise the execution of climate and weather workflows. Exploiting this information automatically may increase the performance, throughput and cost-efficiency of the systems, providing an incentive to users and data-centres that cannot be neglected any longer. Our approach intends to reduce the burden on researchers and, at the same time, optimise the decisions about jobs running on HPC systems.

This paper is structured as follows: First, we describe the software stack involved in executing workflows in climate and weather in Section 1. Next, the vision for including knowledge about data requirements and characteristics is sketched in Section 3 outlining the potential benefit the automatic exploitation may bring. Our design based on existing components in climate and weather is described in Section 4. The paper is concluded in Section 5.

## 1. Workflows in Climate/Weather

In this section, we describe how workflows are executed in a typical software stack and the typical IO stack involved in running a single application.

### 1.1. Cylc

Cylc [18] is in charge of executing and monitoring cyclic workflows in which each step is submitted to the batch scheduler of a data centre. With Cylc, tasks from multiple cycles may be able to run concurrently without violating dependencies preventing the issue of delays that cause one cycle to run into another. Cylc was written in Python and built around a new scheduling algorithm that can manage infinite workflows of cycling tasks without a sequential cycle loop. At any point during workflow execution, only the dependence between the individual tasks matters, regardless of their particular cycle points.

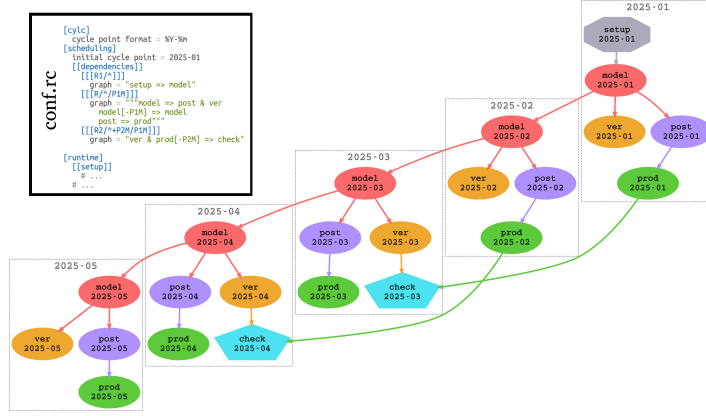
Consider the Cylc workflow for a toy monthly cycling workflow in Fig. 1. In this workflow, an atmospheric model (labelled as **model** in the figure) simulates the physics from a current state to predict the future, for example, a month later. In climate research, this process is repeated in the model to simulate years into the future. Once the simulation of any month is computed, this task completes and implicitly the data for this month becomes available and can now be analysed. In this workflow, the **model** task is followed by postprocessing (**post**), forecast verification (**ver**), and product generation (**prod**) tasks, all specified as a workflow in a Cylc configuration file (`conf.rc`).

The information Cylc uses to control a given workflow is the task dependency. The developers define, in a script file, for each task the parallelism settings and where the data is to be stored.

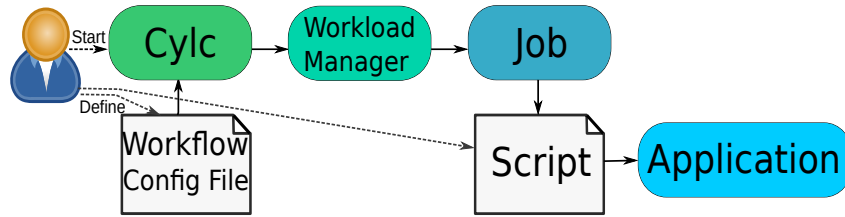
### 1.2. Workflow Execution

While Cylc is directing the execution of workflows, several components are involved in the execution. The software stack involved in executing a general workflow is depicted in Fig. 2. In the following, each stage of the execution is further described.

1. **Scientist** specifies the workflow and provides scripts for each task. After that, the user enacts Cylc to start the workflow.



**Figure 1.** Example of a Cylc workflow with its configuration file

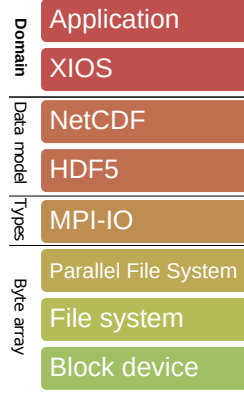


**Figure 2.** Software stack and stages of execution

2. **Cylc** parses the workflow configuration file, generates tasks dependencies, defines a schedule for the execution, and monitors the progress of the workflow. Once a task can be executed (dependencies are fulfilled), the workflow engine submits a *job script* for the workload manager with the required metadata that will run the Cylc task script.
3. **Workload Manager** such as Slurm [9] is responsible to allocate compute resources to a batch job and perform the job scheduling. In our case, it now queues the job that represent the Cylc task and plans its execution considering the scheduling policy of the data centre. Once the job is scheduled to be dispatched, i.e., resources are available, and the job priority is the highest, it is started on the supercomputer.
4. **Job** provides the environment with the resources it runs the user-provided script on one of the nodes allocated for the job. The job sets environment variables containing information about the environment of the batch job, e.g., the compute nodes allocated, and then enacts the Cylc provided script on one node.
5. **Script** executes the shell commands which is one or multiple (potentially parallel) applications to run sequentially. During the creation of the script, Cylc has included variables that describe the task in the workflow. The user script may use commands to create a filename considering the cycle of Cylc or store data in a workflow-specific shared directory. The information is typically fed into the application(s) representing the task, and so defines the storage location.
6. **Application** are executed taking the generated filenames set by the script.

### 1.3. I/O Stack of a Parallel Application

Climate applications may have complex I/O stacks, as can be seen in Fig. 3. In this case, we assume the application uses XIOS [16] which is providing domain-specific semantics to climate



**Figure 3.** I/O path for an MPI-parallel application

and weather. It may gather data from individual fields distributed across the machine (exploiting MPI for parallelism) and then uses NetCDF4 [4] to store data as a file.

Under the hood, NetCDF4 uses the HDF5 API and file format. Internally, HDF5 uses MPI and its data types to specify the nature of the data stored. Finally, data is stored on a parallel file system like Lustre which, on the server-side, stores data in a local file system on block devices such as SSDs and HDDs.

Different applications involved in a workflow may use different I/O stacks to store their outputs. Naturally, the application which uses previously generated data as its inputs must use a compatible API to read the specific data format. In Fig. 3, for example, XIOS may perform parallel I/O via the NetCDF4 API, allowing subsequent processes to read data directly using NetCDF4. Within the ESiWACE project<sup>2</sup>, we are developing the Earth System Data Middleware (ESDM) [13] to allow applications with this kind of software stack to exploit heterogeneous storage resources in a data centre. The goal of ESDM is to provide parallel I/O for parallel applications, with advanced features to optimise subsequent read accesses. Implemented with a standalone API, it also provides NetCDF integration allowing usage in existing applications.

#### 1.4. Data Center Infrastructure

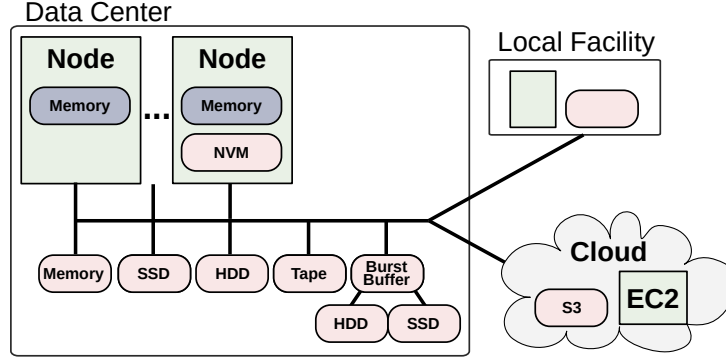
Data centres are providing an infrastructure consisting of computing and storage devices with different characteristics making them more efficient for specific tasks and satisfying the needs for different workflows. Take, for example, the supercomputer Mistral at DKRZ that consists of 3,321 nodes<sup>3</sup>. It offers two types of compute nodes equipped with different CPUs and a range of GPU nodes. Each node provides an SSD for local storage, and DKRZ has additionally two shared Lustre file systems with different performance characteristics. Individual users and projects are mapped to one file system explicitly, and users can access it with work or scratch semantics: While data is kept on the work file system indefinitely, available space is limited by a quota. The scratch file system allows storing more data, but data is automatically purged after some time.

Future centres are expected to have even more heterogeneity. A variety of accelerators (GPU, TPU, FPGAs), active storage, in-memory, and in-network computing technologies will provide storage and processing capabilities. Fig. 4 shows such a system with a focus on computation and storage. Some of these technologies might be local to specific compute nodes or globally

<sup>2</sup><https://www.esiwace.eu/>

<sup>3</sup><https://www.dkrz.de/up/systems/mistral>

available. Depending on the need, the storage characteristics range from predictable low-latency (in-memory storage, NVMe) to online storage (SSD, HDD), and cheap storage for long-term archival (tape). Burst buffer solutions provide a tiered storage system that aims to exploit a storage hierarchy. The tasks within any given workflow could benefit from utilizing different combinations of storage and computing infrastructure.



**Figure 4.** Example of an heterogeneous HPC landscape

## 1.5. Data Management

Usually, the scripts representing tasks define how the data is placed on an available storage system. What happens in many current workflows is that they ignore the benefits of using multiple file systems concurrently and the data locality between tasks to co-locating them. On top of that, in the current state-of-the-art, scientists optimise the available storage resources intuitively and compile the information about this decision-making process manually.

If a user knows the workflow and the system characteristics, s/he can optimise data placement decisions. Consider, for instance, the situation where each computing node has access to three file systems: a fast **scratch** file system on which data may reside only for a week, a slower **work** file system, and a **local** file system. Most current workflows utilize **work** and **scratch** systems. When a task is set to run, the corresponding dataset would be moved from **work** to **scratch**, processed, and the resulting dataset transferred back to **work**. If the **scratch** filesystem reached its capacity, the datasets would be moved back to **work**, and the task would continue running until it is finished, but it might be inefficient. In this situation, there are many obvious opportunities to utilise data migration to optimize performance, and many possible optimization criteria (e.g., costs). However, with a multitude of file systems which differ at each data centre, such optimisations would be difficult to achieve manually for users. Policy-driven systems and burst-buffers perform such optimizations automatically to some extent, however, as they lack information about the workflow, they cannot optimize workflows perfectly.

## 2. State-of-the-Art

Related work can be categorized into: 1) technology that exploits heterogeneous storage environments and supports user-directed policies; 2) solutions for workflow processing.

**Technology.** Manual tiering requires the user or application to control the data placement, i.e., storing data typically in the form of files on a particular storage system and, usually, moving data between storage by scripts. One limitation of such an approach is that decisions about how

data are mapped and packaged into files are made by the producing application, and cannot be changed without manual intervention by a downstream application.

Burst buffer solutions provide a tiered storage system that aims to exploit a storage hierarchy. They can be integrated into hardware capabilities such as DDN IME [3] or pure software solutions. A policy system, e.g., deployed on a burst buffer [21], aims to simplify the data movement for the user, but typically migrates objects in the coarse granularity of files. File systems and data management software such as Spectrum Scale, HPSS, BeeGFS, and Lustre (e.g., using the progressive file layouts feature) provide a hierarchical storage management allowing to store data on different storage technology according to such administrator-provided policies. However, the semantical information that can be used by this type of system to make decisions is limited, e.g., data location, file extension, age of the file, etc.

The storage community had also adjusted various higher-level software to support storage tiering on top of several storage systems, for instance, ADIOS provides in-memory staging that had been exploited by applications [22]. Hermes [11], provides a multi-tiered I/O buffering system with pre-fetcher that provides several data placement policies. iRODS [20], is the rule-oriented data system that allows scientists to organize data into shareable collections and provide metadata such as information about data life cycle. The system provides several patterns for workflows that consider data locality and data migration/replication.

Finally, there have also been extensions to batch schedulers to perform data staging for utilizing node-local storage, for example, NORNS as an extension to Slurm [17].

**Workflows.** A good overview of the flavours of Scientific Workflow Management Systems (SWfMS) and their application to data-intensive workflows is given in [14]. The article states a SWfMS “*should enable the parallel execution of data-intensive scientific workflows and exploit large amounts of distributed resources*”. Existing solutions recognize challenges in data variety (formats of the input data), the opportunity to optimize the schedule by moving code to data, and may consider capacity of available data storage. There have been solutions that allow to specify the data dependencies for tasks. The execution engine Dryad [8], for example, allows to transfer data between tasks via files or directly using TCP connections and attempts to schedule tasks on the same nodes or racks. Swift/T is a scripting language for describing dataflow processing allowing to execute ensembles of applications [19]. Recent improvements aim to migrate data to a local cache allowing to exploit locality, e.g., in [5], information about locality is proposed to be stored in extended attributes. In [15], an approach was presented to monitor and analyze IO behavior of HPC workflows.

Various early research in grid workflows and lately cloud use cases attempts to maximize data locality in that respect. Economic factors (including storage costs) for workflow execution are discussed in [2]. In [6], the authors discuss the role of Machine Learning for workflow execution and elaborate on the general potential for resource provisioning such as optimization of runtime parameters, data movements, and hierarchical storage. An ML model is introduced in [23] that stages data for in-situ analysis by exploiting the access patterns.

Workflow systems can also be utilized specifically for the purpose of reproducing scientific results, i.e., recomputing the results. Those scalable workflow solutions typically utilize a container solution to allow execution in an arbitrary software environment. Popper [10], Snakemake [12], and Nextflow [7] provide a language to specify workflows and to execute them. Snakemake is interesting as it supports to define and infer input and output filenames.

While various aspects of our vision have been addressed individually by related work for different domains, the high level of abstraction that we aim for and the potential it unleashes goes beyond the capabilities of existing approaches.

### 3. Vision for I/O-Aware Workflows

Nowadays, in order to run a job in an HPC environment efficiently, researchers have to develop profound knowledge, not only about their workflow, which is expected, but also about decisions regarding storage, communication, and computing including considerations regarding cost-efficiency of their operations. However, applied scientists should not spend much time understanding hardware characteristics and operative knowledge of running a data centre, but using their time to develop their work and just collect and analyse their results.

We aim for achieving an automatic and dynamic mapping of I/O resources to workflows. Once we have an automated decision about where the job will run and how the storage will be managed, scientists can then reuse their workflow specification on any system without further modification and even without previous knowledge about the system architecture.

There are several approaches to implement the technology for the vision proposed in this work, and changes are needed for software components to realise it. In the next sections, we will discuss a specific design for our transitional roadmap considering climate and weather workflows and tools scientists from this field already use in their routine research.

Our vision for I/O-aware workflows requires two additional pieces of information. Firstly, the user must augment the workflow description with information about I/O requirements and explicitly annotate dependencies to datasets. Secondly, the data centre (or expert user) must provide details about the storage architecture.

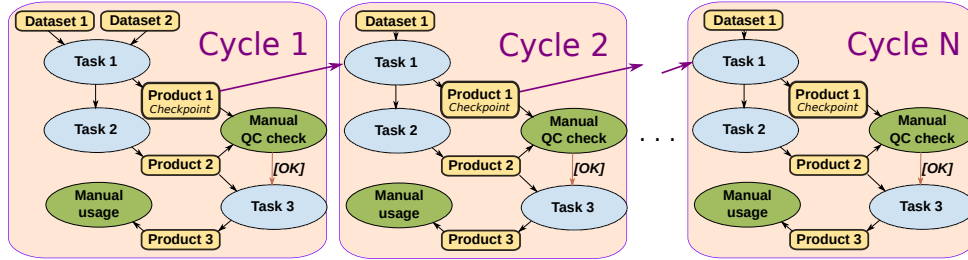
#### 3.1. System Information

While many optimisations are possible once an abstraction is in place, the improvements we discuss here are related to the life cycle of datasets and the placement of such datasets into specific storage according to system performance characteristics and the workflow specification. To achieve that, the system information shall comprise of all available storage systems, the system topology, and details of each of the required components. Simplified and complex models of the components can be included to approximate expected performance for specific I/O patterns. It is expected that the data centre (or expert users) can create such a configuration file, e.g., by using vendor-provided information or by executing benchmarks. With this information, a scheduler can make data placement, transformation, and migration decisions for individual datasets during their life cycle. This separation of concerns allows us to abstract from the workflow what is essential and what a system should optimise to ensure smart usage of available resources.

#### 3.2. Extended Workflow Description

In general, climate and weather workflows allow specifying tasks and the dependencies among them. We aim to enhance the current information with characteristics for input, and output, i.e., the datasets. An example workflow containing input datasets and (intermediate) products is illustrated in Fig. 5. Nodes represent tasks or data, and arrows indicate the dependencies. In the example, Task 1 needs two datasets to perform its work, it directly communicates

with Task 2 and produces Product 1. For each new cycle, the **checkpoint** from the previous cycle, called here **Product 1**, is used as input to starting the next cycle. Most of the workflow can run automatically, except for the manual quality control of the products and the final data usage of Product 3. This last step represents the typical uncertainty of data reuse, i.e., it is unclear how Product 3 will be used further. In the approach proposed in this work, each task is annotated with the required input datasets and the generated products must include metadata such as data life cycle information, the value of data, and how long it should be kept. The idea here is to embrace the concept that tasks dependencies are really imposed by datasets dependencies.



**Figure 5.** Example of a high-level workflow with tasks and data dependencies

### 3.3. Smarter I/O Scheduling

The abstraction and automation of the I/O inside a workflow allows a runtime system to improve data placement and apply data reduction on heterogeneous storage systems. Taking into consideration the architecture and workflow information, a smarter schedule can now be realized by exploiting the additional information. Value and priority can influence fault-tolerance strategies and imply the quality of service for performance and availability. Aspects like data reproducibility (can it be recomputed easily), type of the experiment (test, production), and runtime constraints for the overall and potential workflow could allow reducing costs and, hence, increasing scientific output. Next, we outline the two core strategies and potential the proposed vision can bring to improve current workflows.

**Strategy: Data Placement.** Data placement encompasses all data movement related activities such as transfer, staging, replication, space allocation and de-allocation, registering and unregistering metadata, locating and retrieving data<sup>4</sup>. The general idea is to host a data set on the storage system that is most favorable in terms of performance, cost-effectiveness, availability for the access pattern observed by the workflow. It requires to optimize data locality, where locality is twofold, spatial and temporal on a variety of characteristics. For optimising data placement, there are several approaches:

**Data Allocation** is the assignment of a specific area of an available storage system to particular data. In current workflows, usually, the user has a script for each task defining the filenames with a prefix that places datasets generated by the same task into a specific storage<sup>5</sup>. Because there is one script responsible for generating the configuration, the decision in which directory the dataset will be stored is somewhat fixed. Such configuration is done manually and with restricted information about the system architecture. It would

<sup>4</sup><https://www.igi-global.com/dictionary/data-aware-distributed-batch-scheduling/6782>

<sup>5</sup>Complicated scripts would have allowed changing the storage type depending on the cycle. Still, it is a significant burden to the user.



be interesting to explore storage options for the datasets and, e.g., having datasets from different cycles placed at different storage systems. For instance, in Fig. 5, alternating the storage location for **Product 2** into two scratch file systems is something that would be a simple job for an I/O-aware scheduler. However, currently, that implies having at least two scripts for that task with information about the different storage placement.

**Data Migration** is the process of transferring data from one storage system to another. Data movement involves a large overhead, both in terms of latency and energy-efficient computing, hence needs to be considered carefully. Fig. 6 introduces three possible life cycles for a specific dataset and explains how migrations can be done to improve datasets accessibility. In Fig. 6a, the dataset could be stored on the **local** storage first to avoid congestion on the **work** file system, then be migrated to **work** file system where subsequent tasks of the workflow may read it multiple times. In the end, this dataset might be an intermediate product that can then be deleted. Alternatively (see Fig. 6b), the dataset could be stored on the **scratch** file system immediately and accessed there. However, that would require subsequent tasks to be placed on the same node where data is now stored, and the last read access must happen before files on **scratch** are automatically removed. The last scenario (see Fig. 6c) presents the case where the dataset is created on **work** by a task and it is copied to a **local** node. This **local** node allows reading accesses of subsequent tasks which might be beneficial for small random accesses. However, it would require that the subsequent tasks are placed on the same node where data is now stored.

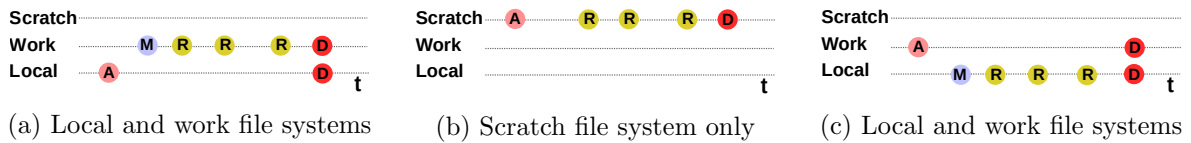
**Data Replication** in computing involves sharing information to ensure consistency between redundant resources, such as software or hardware components, to improve reliability, fault-tolerance, or accessibility. Data might be replicated by enabling the system to rerun parts of the workflow in case of a data loss. In addition, the system may combine the replication of data by **transforming** the data into a different representation allowing to achieve better performance for a variety of access patterns.

**Direct-Coupling** replaces I/O by communicating data between subsequent steps of a workflow directly without storing intermediate data products on persistent storage. As an example, in Fig. 5, the outcome of Task 1 may be used directly by Task 2. To achieve a decoupling between producer and consumer, it may also be kept in memory and cached.

In the design proposed in this work (Section 4), we will focus on the data placement strategy.

**Strategy: Data Reduction.** Data reduction reduces the amount of data stored. We discuss here two potential optimisations: **data compression**, and **data recomputation**.

**Data Compression** is the process of encoding information using fewer bits than the original representation. Knowing the characteristics of data production and usage makes it simpler



**Figure 6.** Alternative life cycles for mapping a dataset to storage and the operations: **A**llocation, **M**igration, **R**eading, and **D**eleting

for scientists to annotate the required precision of data in those workflows. The storage system can exploit such information by reducing precision of data and automatically picking an appropriate compression algorithm.

**Data Recomputation** Climate/weather scientists are trading recomputation with space usage manually. By knowing how to rerun the workflow behind the data creation, a smarter storage system can automatically trade data availability for potential recomputation opportunities to optimise the cost-efficiency of the system. Intermediate states could be rerun by utilising virtualisation and container technologies. Consider Fig. 5 again and that, at every  $K$  cycles of the workflow, the generated **Product 3** (from Cycle 1 to Cycle  $K$ ) are used in a validation task, called here **check**. From the workflow<sup>6</sup>, we know that  $P_3C_1$  will be used to construct  $P_3C_2$  and then **check**. This dataset will probably be stored somewhere, and it will not be used until the workflow reaches the  $K$ -th cycle. One alternative is to delete it after it was first used and then recompute it when time is right. The cost of doing that is storing **checkpoint** and then use it to reconstruct product  $P_3C_1$ . If, for instance,  $P_3C_1$  is a large dataset, **checkpoint** is small, and computing time is short. It is easy to see that deleting and recomputing it may improve the costs for running a workflow. This is just an example, and, currently, scientists perform optimisations manually.

### 3.4. Benefit

The benefits of the proposed vision are:

**Abstraction** The user does not have to know the architecture of the target system on which the workflow will run, removing the specialist from the decision-making process.

**Optimisation** The workflow will be optimised specifically for the available system infrastructure and extra information about the data. In particular, by exposing the heterogeneous architecture, potentially runtime characteristics can be considered. By using information about the value of data, policies for data management (storage resilience, recomputation, replication, etc.) can be decided.

**Performance-portability** With both abstraction and optimisation, the user can specify the I/O requirements only once for the tasks of a specific workflow, and the workflow can now run on any system without user intervention. Even more, if the system characteristics change, e.g., it gets upgraded, an additional storage tier becomes available, or if storage degrades, the I/O-aware workflow could automatically adapt and make use of this new environment.

## 4. Design

This section describes our first approach to incrementally extend workflows for climate and weather that realises parts of our vision. While individual components such as ESDM and Cylc exist, we have not implemented the described scheduler, yet. To automatically make scheduling decisions, the software stack needs to:

1. Deliver information about dataset life cycle together with the workflow, and

---

<sup>6</sup>The  $P_iC_j$  notation represents the Product  $i$  generated in the Cycle  $j$ .

2. Adapt the resulting workflow, individual scripts and application executions to consider the potential for data placement and migration.

#### 4.1. System Information

ESDM is used as IO middleware in the parallel application (in NetCDF or directly) and orchestrates the IO according to a configuration file. The information about the system is provided using a simplified ESDM configuration file (`esdm.conf`) provided by the data centre (or expert users). This file contains information about the available technology in the data centre, its I/O characteristics, and will be used to make decisions about how to prioritise I/O targets.

In the example presented in Listing 1, we have three storage targets: two global accessible file systems (`lustre01` and `lustre02`), and one local file system in `/tmp` that can be accessed via the POSIX backend. Each of them comes with a lightweight performance model and the maximum size of data fragments. The metadata section (Line 24) utilises here a POSIX interface to store the information about the ESDM objects. Internally, ESDM creates so-called containers and dataset objects to manage data fragments.

**Listing 1.** Example of an ESDM configuration file (`esdm.conf`)

```

1  "backends": [
2      {"type": "POSIX", "id": "work1", "target": "/work/lustre01/projectX/",
3         "performance-model" : {"latency" : 0.00001, "throughput" : 500000.0},
4         "max-threads-per-node" : 8,
5         "max-fragment-size" : 104857600,
6         "max-global-threads" : 200,
7         "accessibility" : "global"
8     },
9     {"type": "POSIX", "id": "work2", "target": "/work/lustre02/projectX/",
10        "performance-model" : {"latency" : 0.00001, "throughput" : 200000.0},
11        "max-threads-per-node" : 8,
12        "max-fragment-size" : 104857600,
13        "max-global-threads" : 200,
14        "accessibility" : "global"
15    },
16    {"type": "POSIX", "id": "tmp", "target": "/tmp/esdm/",
17        "performance-model" : {"latency" : 0.00001, "throughput" : 200.0},
18        "max-threads-per-node" : 0,
19        "max-fragment-size" : 10485760,
20        "max-global-threads" : 0,
21        "accessibility" : "local"
22    }
23 ],
24 "metadata": {"type": "POSIX",
25     "id": "md",
26     "target": "./metadata",
27     "accessibility" : "global"
28 }
```

ESDM manages a pool of threads that should be created per compute node to achieve good performance and delegates the assignment of optimal block sizes to the storage backend. The number of threads is defined in the configuration file. As an example, based on the number of Object Storage Targets (OSTs) available to both Lustre systems at DKRZ, performance tests already developed by ESDM [1] shows that no more than 200 threads in total should be used

to perform I/O to extract the best performance. To clarify the behaviour, ESDM distributes a single dataset across multiple storage devices depending on their characteristics. Since ESDM also supports several (non-POSIX) storage backends, an application can utilise all available storage systems without any modifications to the code.

The configuration file is inquired by an application that utilises ESDM and steers the distribution of data during I/O. While the current system information and performance model is based on latency and throughput only, it shows that automatic decision making can be made on behalf of the user.

## 4.2. Extended Workflow Description

The user now has to provide information about datasets required for input and the generated output for each Cylc task in a separate file similarly to Cylc's workflow configuration file. An example of an I/O-workflow configuration file is shown in Listing 2. In this file, information about Task 1 is given by example, and we expect the extra information about all tasks in the same file; ultimately, this could be integrated into the workflow specification file of Cylc.

**Listing 2.** External Cylc I/O-workflow configuration file

```

1 [Task 1]
2   [[inputs]]
3     topography = "/pool/input/app/config/topography.dat"
4     checkpoint = "[Task 1].checkpoint$(CYCLE - 1)"
5     init = "/pool/input/app/config/init.dat"
6
7   [[outputs]]
8     [[varA]] # This is the name of the variable
9       pattern = 1 day
10      lifetime = 5 years
11      type = product
12      datatype = float
13      size = 100 GB
14      precision.absolute_tolerance = 0.1
15
16     [[[checkpoint]]]
17       pattern = $(CYCLE)
18       lifetime = 7 days
19       type = checkpoint
20       datatype = float
21       dimension = (100,100,100,50)
22
23     [[[log]]]
24       type = logfile
25       datatype = text
26       size = small

```

In this example, the workflow file could define a cycle flexibly to be a month or a year according to the Cylc workflow file. The notation is similar to the specification of Cylc workflows using a nested INI format. For each task, inputs and outputs are defined. In the input section, each entry specifies the virtual name that is used by ESDM as a filename inside NetCDF. Line 3, for example, defines that the filename **topography** is mapped to a specific input file. This dataset does not depend on any previous step of the workflow. The next line specifies that the input

filename **checkpoint** should be mapped to the output of Task 1 checkpoint dataset from the previous cycle (e.g., the checkpoint generated after completed the last year’s output). For the initial cycle, the checkpoint file will be empty, and the application will load the init data. In the output section, the datasets are annotated with their characteristics more precisely. For each variable, a pattern defining how frequently the data is output according to the workflow must be provided. Most data is input and output in the periodicity of the cycle, except for **varA**, which is output per day regardless of the cycle. Next, we formally define the expected annotations in all the fields expected in the I/O-workflow configuration file:

**Name** A basic name for the field/data generated (Lines: 8, 16, 23). It is extended by a pattern defined in a variable.

**Pattern** The frequency the data is output (Lines: 9, 17).

**Lifetime** How long the data must be retained on storage (if at all) (Lines: 10, 18).

**Type** The class type of data, i.e., checkpoint, diagnostics, temporary (Lines: 11, 19, 24).

**Datatype** The data type of the data (Lines: 12, 20, 25).

**Size** An estimate of the data size (Lines: 13, 26). This field can be inferred if dimension and datatype are provided.

**Dimension** The data dimension (Line 21).

**Accuracy** This includes various characteristics quantifying the required level of data precision; For example, in Line 14.

Note that the user may not be able to provide all required information. This can be handled by assuming a default safe behavior. For instance, in case of a missing data precision, data should be retained in original form. Knowing the dimension or size a priori might be difficult for scientists, e.g., the log file size is unclear. In this case, the user may insert relative information like small or big, reflecting that any information is better than no information at all. In future, we will explore how to automatically infer the output volume from the input or by using monitoring. By allowing to run using an empty I/O workflow specification and monitoring I/O accesses for one cycle, we can propose an I/O description to the user to simplify the specification.

### 4.3. Smarter I/O Scheduling

From the list of opportunities, we realise data placement and migration in a heterogeneous (multi-storage) environment. These goals will be achieved via the proposed I/O-aware scheduler, called here EIOS (ESDM I/O Scheduler). EIOS will make the schedule considering Cylc workflow, and ESDM provided system characteristics. Components of it are involved in different steps of the workflow and the I/O path.

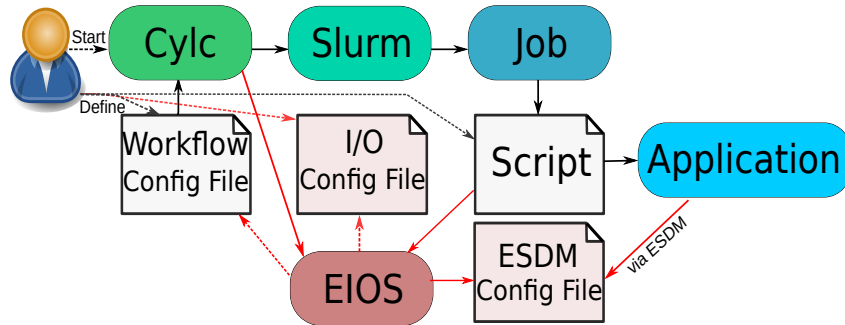
While Cylc schedules the workflow, EIOS can provide hints about co-locating tasks which provide the opportunity for keeping data in local storage. Decisions about data locality will not be made for a whole (and potentially big) workflow. Instead, the system will make decisions by looking ahead to several steps of the workflow, allowing reacting to the observed dynamics of the execution. We utilise DDN’s IME API to pin data in IME and to trigger migrations between IME and a storage backend explicitly.

Ultimately, when a user-script runs, the information about the intended I/O schedule is communicated by modifying the "filename" which is then used by the ESDM-aware application to determine the I/O scheduler.

#### 4.4. Modified Workflow Execution

The steps to execute a workflow enriched with IO information and performing smarter scheduling is depicted in Fig. 7. The suggested alterations can be seen in boxes pointed by red arrows, and the remaining components are the current state-of-the-art for workflows in climate and weather from Fig. 2.

In the following, we describe the modifications we propose in this vision paper for each component involved in the software stack.



**Figure 7.** Software stack and stages of execution with our I/O-aware scheduler – EIOS

1. **Scientist** The end-user now has to provide an additional file that covers the I/O information for each task and slightly changes (discussed above) have to be made to the current scripts.
2. **Cylc** EIOS is invoked by Cylc to identify potential optimisations in the schedule before generating the Slurm script.
3. **EIOS** The ESDM I/O Scheduler reads the information about the workflow (original and I/O workflow configuration files) and acts depending on the stage of the execution. EIOS consists of several subcomponents:

- The high-level scheduler that interfaces with Cylc.
- A tool to generate pseudo-file names used by the ESDM-aware applications.
- A data management service (not shown on the figure) that migrate and purge data at the end of the life cycle.

EIOS components use knowledge about the system by parsing the `esdm.conf` file. EIOS may decide that subsequent jobs shall be placed on the same node, reorder the execution of some jobs, and provide information for conducting data migration.

4. **Slurm** Cylc may now have added an optimisation identified by EIOS, and that is now handled by a modified Slurm. Also, if migrations have to be performed, Slurm will administer them according to the specification in the job script.
5. **Job** A job might run on the same nodes than a previous job to utilise local storage.
6. **Script** Filenames are now generated by a replacement command that calls EIOS to create a pseudo filename. This filename will encode additional information for ESDM about how to prioritise data placement according to data access.
7. **Application** The application may either use XIOS, NetCDF with ESDM support or ESDM directly to access datasets. Hence, in Fig. 3, the HDF5 layer is replaced with ESDM. ESDM

loads the `esdm.conf` file that contains the information about the available storage backends and their characteristics. ESDM extracts the long-term schedule information from the generated pseudo filenames and employs it during the I/O scheduling to optimise the storage considering data locality between tasks.

## 5. Conclusions

Organising the data placement on storage tiers is particularly in the domain of climate and weather performed by the users or via policies, leading often to suboptimal decisions. Additionally, the manual optimization and hard-coding of storage locations is non-portable and an error-prone task. We believe users must be able to express their workflows in an abstract fashion. By increasing the abstraction level for scientists, not only tedious manual optimisation could be automatised, but also strategies for data placement and data reduction can be harnessed. With knowledge about the data pattern, the runtime system could generate optimized execution plans and monitor their execution. In this work, we describe the general vision and a specific design for the software stack in the domain in climate and weather that we work on in the ESiWACE project. The changes proposed increase the opportunity for smarter scheduling of storage in heterogeneous storage environments by considering the characteristics of data and system in a workflow.

## Acknowledgements

This project is funded by the European Union’s Horizon 2020 research and innovation programme under grant agreement No. 823988. We thank our collaborator Bryan Lawrence for their input to this article.

*This paper is distributed under the terms of the Creative Commons Attribution-Non Commercial 3.0 License which permits non-commercial use, reproduction and distribution of the work without further permission provided the original work is properly cited.*

## References

1. Final Implementation of the Earth System Data Middleware (D4.3) (Jun 2019), DOI: 10.5281/zenodo.3361225, ESiWACE has received funding from the European Union’s Horizon 2020 Research and Innovation Programme under Grant Agreement No 675191
2. Alkhanak, E.N., Lee, S.P., Rezaei, R., Parizi, R.M.: Cost optimization approaches for scientific workflow scheduling in cloud and grid computing: A review, classifications, and open issues. *Journal of Systems and Software* 113, 1–26 (2016)
3. Betke, E., Kunkel, J.: Benefit of DDN’s IME-Fuse and IME-Lustre File Systems for I/O Intensive HPC Applications. In: Yokota, R., Weiland, M., Shalf, J., Alam, S. (eds.) *High Performance Computing: ISC High Performance 2018 International Workshops*, Frankfurt/Main, Germany, June 28, 2018, Revised Selected Papers. pp. 131–144. No. 11203 in *Lecture Notes in Computer Science*, ISC Team, Springer (01 2019), DOI: [https://doi.org/10.1007/978-3-030-02465-9\\_9](https://doi.org/10.1007/978-3-030-02465-9_9)
4. Center, U.P.: Network common data form (netcdf), DOI: <http://doi.org/10.5065/D6H70CW6>
5. Dai, D., Ross, R., Khaldi, D., Yan, Y., Dorier, M., Tavakoli, N., Chen, Y.: A Cross-Layer Solution in Scientific Workflow System for Tackling Data Movement Challenge. *arXiv preprint arXiv:1805.06167* (2018)

6. Deelman, E., Mandal, A., Jiang, M., Sakellariou, R.: The role of machine learning in scientific workflows. *The International Journal of High Performance Computing Applications* 33(6), 1128–1139 (2019)
7. Di Tommaso, P., Chatzou, M., Floden, E.W., Barja, P., Palumbo, E., Notredame, C.: Nextflow enables reproducible computational workflows. *Nature Biotechnology* 35, 316–319 (04 2017), DOI: 10.1038/nbt.3820
8. Isard, M., Budiu, M., Yu, Y., Birrell, A., Fetterly, D.: Dryad: distributed data-parallel programs from sequential building blocks. In: *Proceedings of the 2nd ACM SIGOPS/EuroSys European Conference on Computer Systems* 2007. pp. 59–72 (2007)
9. Jette, M.A., Yoo, A.B., Grondona, M.: Slurm: Simple linux utility for resource management. In: *Lecture Notes in Computer Science: Proceedings of Job Scheduling Strategies for Parallel Processing (JSSPP)* 2003. pp. 44–60. Springer-Verlag (2002)
10. Jimenez, I., Sevilla, M., Watkins, N., Maltzahn, C., Lofstead, J., Mohror, K., Arpaci-Dusseau, A., Arpaci-Dusseau, R.: The popper convention: Making reproducible systems evaluation practical. In: *Parallel and Distributed Processing Symposium Workshops (IPDPSW)*, 2017 IEEE International. pp. 1561–1570. IEEE (2017)
11. Kougkas, A., Devarajan, H., Sun, X.H.: I/o acceleration via multi-tiered data buffering and prefetching. *Journal of Computer Science and Technology* 35(1), 92–120 (2020)
12. Köster, J., Rahmann, S.: Snakemake: A scalable bioinformatics workflow engine. *Bioinformatics* 28(19), 2520–2522 (08 2012), DOI: 10.1093/bioinformatics/bts480
13. Lawrence, B.N., Kunkel, J.M., Churchill, J., Massey, N., Kershaw, P., Pritchard, M.: Beating data bottlenecks in weather and climate science. In: *Extreme Data Workshop – Forschungszentrum Jülich, Proceedings, IAS series, volume 40*. pp. 31–36 (2018)
14. Liu, J., Pacitti, E., Valduriez, P., Mattoso, M.: A survey of data-intensive scientific workflow management. *Journal of Grid Computing* 13(4), 457–493 (2015)
15. Lüttgau, J., Snyder, S., Carns, P., Wozniak, J.M., Kunkel, J., Ludwig, T.: Toward Understanding I/O Behavior in HPC Workflows. In: *IEEE/ACM 3rd International Workshop on Parallel Data Storage & Data Intensive Scalable Computing Systems (PDSW-DISCS)*. pp. 64–75. IEEE Computer Society, Washington, DC, USA (02 2019), DOI: <https://doi.org/10.1109/PDSW-DISCS.2018.00012>
16. Meurdesoif, Y., Caubel, A., Lacroix, R., D’erouillat, J., Nguyen, M.H.: Xios tutorial (2016), <http://forge.ipsl.jussieu.fr/ioserver/raw-attachment/wiki/WikiStart/XIOS-tutorial.pdf>
17. Miranda, A., Jackson, A., Tocci, T., Panourgias, I., Nou, R.: Norns: Extending slurm to support data-driven workflows through asynchronous data staging. In: *2019 IEEE International Conference on Cluster Computing (CLUSTER)*. pp. 1–12. IEEE (2019)
18. Oliver, H., Shin, M., Matthews, D., Sanders, O., Bartholomew, S., Clark, A., Fitzpatrick, B., van Haren, R., Hut, R., Drost, N.: Workflow automation for cycling systems: The cylc workflow engine. *Computing in Science Engineering* 21(4), 7–21 (July 2019), DOI: 10.1109/MCSE.2019.2906593
19. Ozik, J., Collier, N.T., Wozniak, J.M., Spagnuolo, C.: From desktop to large-scale model exploration with swift/t. In: *2016 Winter Simulation Conference (WSC)*. pp. 206–220. IEEE (2016)
20. Rajasekar, A., Moore, R., Hou, C.y., Lee, C.A., Marciano, R., de Torcy, A., Wan, M., Schroeder, W., Chen, S.Y., Gilbert, L., et al.: iRODS primer: integrated rule-oriented data system. *Synthesis Lectures on Information Concepts, Retrieval, and Services* 2(1), 1–143 (2010)



21. Romanus, M., Ross, R.B., Parashar, M.: Challenges and considerations for utilizing burst buffers in high-performance computing. CoRR abs/1509.05492 (2015), <http://arxiv.org/abs/1509.05492>
22. Slawinska, M., Clark, M., Wolf, M., Bode, T., Zou, H., Laguna, P., Logan, J., Kinsey, M., Klasky, S.: A Maya use case: adaptable scientific workflows with ADIOS for general relativistic astrophysics. In: Proceedings of the Conference on Extreme Science and Engineering Discovery Environment: Gateway to Discovery. pp. 1–8 (2013)
23. Subedi, P., Davis, P.E., Parashar, M.: Leveraging machine learning for anticipatory data delivery in extreme scale in-situ workflows. In: 2019 IEEE International Conference on Cluster Computing (CLUSTER). pp. 1–11. IEEE (2019)